

---

# Combining Density and Overlap (CoDO): A New Method for Assessing the Significance of Overlap Among Subgraphs

---

**Abram Magner**<sup>\*,†</sup>  
Coordinated Science Lab, UIUC  
Champaign, IL 61820  
anmagner@illinois.edu

**Shahin Mohammadi**<sup>†</sup>  
Dept. of Computer Science, Purdue University  
West Lafayette, IN 47907  
mohammadi@purdue.edu

**Ananth Grama**  
Dept. of Computer Science, Purdue University  
West Lafayette, IN 47907  
ayg@cs.purdue.edu

## Abstract

Algorithms for detecting clusters (including overlapping clusters) in graphs have received significant attention in the research community. A closely related important aspect of the problem – quantification of statistical significance of overlap of clusters, remains relatively unexplored. This paper presents the first theoretical and practical results on quantifying statistically significant interactions between clusters in networks. Such problems commonly arise in diverse applications, ranging from social network analysis to systems biology. The paper addresses the problem of quantifying the statistical significance of the observed overlap of the two clusters in an Erdős-Rényi graph model. The analytical framework presented in the paper assigns a  $p$ -value to overlapping subgraphs by combining information about both the sizes of the subgraphs and their edge densities in comparison to the corresponding values for their overlapping component. This  $p$ -value is demonstrated to have excellent discrimination properties in real applications and is shown to be robust across broad parameter ranges.

Our results are comprehensively validated on synthetic, social, and biological networks. We show that our framework: (i) derives insight from both the density and the size of overlap among communities (circles/pathways), (ii) consistently outperforms state-of-the-art methods over all tested datasets, and (iii) when compared to other measures, has much broader application scope. In the context of social networks, we identify highly interdependent (social) circles, and show that our predictions are highly co-enriched with known user features. In networks of biomolecular interactions, we show that our method identifies novel cross-talk between pathways, sheds light on their mechanisms of interaction, and provides new opportunities for investigations of biomolecular interactions.

## 1 Introduction

Quantification of statistical significance of an observed artifact in data is a critical aspect of data analytics applications, particularly at scale. The statistical significance of an artifact is often quantified

---

<sup>\*</sup>Corresponding authors: anmagner@illinois.edu

<sup>†</sup>Authors contributed equally.

in terms of its  $p$ -value – the probability that a random event caused the observed artifact. Clearly, the lower this probability, the higher the statistical significance. Statistical significance, in terms of its  $p$ -value is defined with respect to a prior, which models the background (random) distribution. Selecting a prior that is distant from the data renders all observations significant (i.e., they have very low  $p$  values). Conversely, selecting a prior that is identical to the data renders all observations insignificant (i.e., their  $p$  values approach 1). A suitable combination of a statistical prior along with an appropriate  $p$  value formulation provides discrimination of significant artifacts from those that are not.

$p$ -values of observed artifacts may be estimated using analytic or sampling techniques. Sampling techniques for artifacts in graphs often require large numbers of samples for convergence. Analytic techniques, on the other hand, pose significant challenges for derivation of tight bounds, particularly for complex graph models. Analytic techniques have been successfully demonstrated in the context of characterizing the significance of dense subgraphs [1], distributions of subgraph diameters, and frequencies of network motifs such as triangles. Sampling based techniques have been demonstrated in conjunction with non-parametric graph priors to characterize the alignment of networks [2].

In this paper, we focus on the important and challenging problem of characterizing the statistical significance, in terms of its  $p$ -value, of the observed overlap between two dense subgraphs. This problem arises in a number of applications – for instance: (i) in understanding when an overlap between two communities in a social network suggests a social tie; (ii) in assessing whether overlap among two groups of like-minded movie viewers suggests a shared interest; and (iii) in determining whether crosstalk between two sets of interacting biomolecules corresponds to a functional property.

The problem of analytically characterizing a suitable  $p$ -value for the overlap of two subgraphs, even for simple graph models such as Erdős-Rényi, is a hard one. One can use simpler abstractions to characterize observed overlap using Hypergeometric Tails (HGT) or by characterizing the significance of the overlapping subgraph independently (i.e., ignoring the existence of the two sub-graphs). However, these simpler abstractions lack the discriminating power necessary for most applications. The Hypergeometric Tail does not consider network structure, therefore its applicability is limited. Characterizing the significance of the overlapping subgraph does not consider cluster sizes, relying, instead on high density within the overlapping subgraph. Our method, on the other hand, allows for cases in which densities of each component subgraph and their overlap are not significant by themselves with respect to the ambient graph, but the edges inside the two subgraphs are unusually concentrated inside the overlap. This aspect of our method gives it excellent robustness across wide parameter ranges, as well as diverse applications. We present a detailed theoretical foundation for quantifying the  $p$ -value of the overlap of two subgraphs in an Erdős-Rényi graph. Our formulation considers the sizes of the two sub-graphs, their densities, as well as the size and density of the overlap set. In doing so, it provides excellent discrimination across the parameter space.

We comprehensively validate the suitability of our selected prior and our framework for quantifying statistical significance of overlaps on synthetic as well as real networks. Using an Erdős-Rényi random graph with implanted clusters of different size, density, and overlap, we experimentally evaluate the behavior of our framework with respect to different parameters. Having established the baseline behavior, we apply our method to the ego networks for Facebook, Google+, and Twitter networks [3]. We identify social circles with significant overlap and show that these circles are highly enriched in common features. Finally, using a well curated network of interacting human proteins and a given set of functional pathways, we identify statistically significant crosstalk between these pathways. We show that the crosstalk identified correlates very well with known biological insights, while also providing a number of novel observations for future investigation. We also demonstrate that our measure significantly outperforms other characterizations of statistical significance.

## 2 Background and Related Work

Evaluating the significance of overlap among a pair of given subgraphs is a challenging problem that can be addressed from different viewpoints, depending on how significance is quantified. The most common method for assessing overlap is to ignore the underlying interactions and focus on the size of the overlap (number of vertices), compared to the sizes of subgraph pairs. In this approach, one often uses Fisher’s exact test, also known as the hypergeometric test [4]. This method is described in more detail in Section 3.1. Another way to define significance of overlap is in terms of density. Here,

an overlap is defined to be significant if there are “many” edges in the overlap set, compared to a null model. In this class of methods the size of overlap set is assumed to be fixed and density of overlap is assessed, conditioned on the fixed size. The key aspect of this measure is the the density of overlap.

There has been extensive research focused on the algorithmic task of identifying dense components in graphs (see, e.g., [5] for a survey). Regarding the study of dense subgraphs in random graphs, the papers by Arratia in 1990 and Koyutürk et al. in 2007 [1, 6] are perhaps the most relevant. Both study the typical sizes of subgraphs of a given density in Erdős-Rényi graphs, motivated by the study of protein interaction networks. The former also gives a result on overlaps between dense subgraphs, but only in the context of dense graphs, unlike social or biological networks that are sparse and the number of edges and nodes are typically of the same order. The distributions of the maximum sizes of cliques and independent sets in random graphs have also been studied, essentially completely [7]. The typical behavior of the overlaps between independent sets in the case of sparse random graphs (equivalently, cliques in extremely dense random graphs) has recently been studied [8]. However, the range of parameters for which the results are derived differs significantly from ours, rendering these results less relevant to our applications of interest. While these contributions are fundamental to our understanding of statistics regarding the distribution of dense components in a graph, *none* of them directly address the problem of assessing the overlap among subgraphs. Here, we combine the two approaches based on sizes of the overlapping subgraphs and the density of the overlap compared to the constituent subgraphs to analytically derive a  $p$ -value for the overlap that considers both its size and density.

From an applications point of view, this problem has important implications in social networks, systems biology, financial transactions, and network security. In systems biology, as pathways associated with specific biological functions are fully resolved, there is increasing need to assess the extent of overlap/ cross-talk among these pathways. Pathway interactions have been shown to play a key role in development and progression of cancers [9, 10]. Different groups have focused on identifying the cross-talk map among pathways [11, 12]. However, to the best of our knowledge, none of these methods directly use the overlap subgraph to infer pathway interactions. Another important application of our method in life sciences is in geneset enrichment analysis. While the hypergeometric test remains the most widely used method for identifying functional enrichment of a set of differentially expressed genes, most recent methods do not focus only on the geneset overlap but also try to incorporate network context to evaluate functional importance of genesets [13, 14]. However, a majority of these methods rely on computing empirical  $p$ -values instead of providing a closed-form solution. Our method is the first method that quantifies statistical significance within a network context, and in doing so, provides significantly higher discriminating power than prior methods, while using significantly fewer computing resources.

### 3 Statistical significance of subgraph overlaps

Given a pair of overlapping dense subgraphs in a graph, our goal is to derive a measure of the extent to which we should be surprised by their overlap. The appropriate framework for this is that of *statistical significance*. In the most general setting, we have a random variable  $X$  on a probability space and an observed value  $\hat{X}$  for  $X$ , and we want a measure of the surprise in observing  $\hat{X}$ . To the observed value  $\hat{X}$ , we associate a  $p$ -value, which is given by the probability that  $X$  takes a value *as extreme or more extreme* than  $\hat{X}$  (the notion of extremity depends on the range of the random variable). If the  $p$ -value is less than some fixed threshold, the observed value is said to be *statistically significant* with respect to the distribution of  $X$  (the *prior*); that is, the observation  $\hat{X}$  is unlikely to have occurred by chance.

For a given random variable (which need not even take values in a partially ordered set), there may be many non-equivalent notions of  $p$ -value that take into account various pieces of observable information about  $X$ . In general, a notion of statistical significance is preferable to another if it has more *discriminating power*, in the sense that the  $p$ -value decays smoothly as observations become more extreme (this allows for statistically significant observations to be detected *and* to be ranked by relative significance).

In what follows, we will describe three formulations of statistical significance for subgraph overlaps. The first two, which are already present in literature, take into account only partial information: the first takes into account only the sizes of the subgraphs and their intersection; the second takes into

account only the size and density of the intersection of the overlapping subgraphs. The third, which we call *CoDO* and which is the main topic of this paper, combines subgraph sizes and densities of individual subgraphs along with their overlaps, to yield a test with significantly greater discriminating power.

The basic setup in all three tests is as follows: we have a graph  $G \sim \mathcal{G}(n, p)$  (where  $\mathcal{G}(n, p)$  denotes the Erdős-Rényi distribution on graphs of size  $n$ , with edge probability  $p$ ) and subgraphs  $A$  and  $B$  with  $X = A \setminus B$ ,  $Y = B \setminus A$ , and  $Z = A \cap B$ . We also have the following definitions and notation relating to density: For any subgraph  $S \subseteq G$ , we denote by  $E(S)$  the set of edges in  $S$ , and  $e(S) = |E(S)|$ . For a pair of subgraphs  $S_1, S_2$ , we denote the set of edges between nodes in  $S_1$  and nodes in  $S_2$  by  $E(S_1, S_2)$ , and  $e(S_1, S_2) = |E(S_1, S_2)|$ . We denote by  $\delta(S)$  the *density* of  $S$ :

$$\delta(S) = \frac{e(S)}{\binom{|S|}{2}}.$$

Similarly, for  $S_1, S_2$ ,

$$\delta(S_1, S_2) = \frac{e(S_1, S_2)}{|S_1||S_2|}.$$

We now describe several  $p$ -value formulations.

### 3.1 Hypergeometric Tail (HGT)

The *hypergeometric tail*  $p$ -value takes into account only the *size* of the overlap  $Z$  of  $A$  and  $B$ . It is defined by considering the following probabilistic experiment (note that this hypothetical experiment is introduced for the purpose of defining a particular probability distribution as the outcome of some random process; it is *not* meant to be implemented by a user of our methods): we fix a subset  $\hat{B} \subseteq V(G)$  of size  $|\hat{B}| = |B|$ , and we draw, uniformly at random from  $V(G)$ , a subset  $\hat{A}$  of nodes of size  $|\hat{A}| = |A|$ . The size  $|\hat{A} \cap \hat{B}|$  of the overlap is then hypergeometrically distributed, so the hypergeometric  $p$ -value is given by the probability that the overlap in this experiment has size at least that of the observed overlap  $Z$ . That is,

$$\Pr[|\hat{A} \cap \hat{B}| \geq |Z|] = 1 - F(|Z| - 1, n, |\hat{B}|, |\hat{A}|),$$

where  $F(x, y, z, w)$ , for arbitrary  $x \in \mathbb{N} \cup \{0\}$ , denotes the probability that a hypergeometrically distributed random variable with population size  $y$ , number of successes  $z$ , and number of trials  $w$ , takes a value at least  $x$ .

### 3.2 Erdős-Rényi Density Model (ERD)

A second approach to scoring of  $p$ -values of overlaps takes into account only the size and density of the overlap set  $Z$ . It is defined as the probability that there exists in  $G$  a subgraph of density  $\delta(Z)$  and size at least  $|Z|$ . This is given by

$$\Pr[\exists H \subseteq G, |H| \geq |Z| : \delta(H) = \delta(Z)].$$

Upper bounds on this probability have been worked out by, e.g., Arratia in 1990 and Koyutürk et al. in 2007 [1, 6]. The main tool for this is the *first moment method* [15]: the event that there exists a subgraph with a given property is precisely equal to the event that the number of subgraphs with the given property is at least 1. The first moment method consists of an application of Markov's inequality to conclude that

$$\Pr[\#\{H : |H| \geq |Z|, \delta(H) = \delta(Z)\} \geq 1] \leq \mathbb{E}[\#\{H : |H| \geq |Z|, \delta(H) = \delta(Z)\}]$$

The expected value is then easily calculated using linearity of expectation. For a fixed value of  $\delta(Z)$ , the threshold value of  $|Z|$  below which the expected value tends to  $\infty$  as  $n \rightarrow \infty$  and above which it tends to 0 turns out to be  $\frac{2 \log(n)}{\kappa(\delta(Z), p)}$ , where  $\kappa(a, b)$  is given by

$$\kappa(a, b) = a \log(a/b) + (1 - a) \log((1 - a)/(1 - b)), \quad (1)$$

the relative entropy between Bernoulli random variables with parameter  $a$  and  $b$ , respectively. The first moment method shows that overlaps with density  $\delta(Z)$  and size at least  $(1 + \epsilon)|Z|$ , for any

fixed positive  $\epsilon$ , are unlikely to occur randomly. The *second moment method*, which takes into account the dependence between events defined on overlapping subgraphs, can be used to give a matching probabilistic lower bound, which shows that, with high probability (i.e., with probability asymptotically tending to 1), the maximum size of any  $\delta(Z)$ -dense subgraph is asymptotically equivalent to  $\frac{2 \log(n)}{\kappa(\delta(Z), p)}$ .

Because of the very sharp transition, as  $|Z|$  increases, from insignificant to significant, this formulation has the undesirable property that it has little discriminating power. In particular, the first moment upper bound implies that

$$\Pr[\#\{H : |H| \geq |Z|, \delta(H) = \delta(Z)\} \geq 1] \leq e^{-\Theta(|Z|^2)}.$$

That is, the  $p$ -value decays superexponentially as  $|Z|$  increases.

### 3.3 Combining Density/Overlap (CoDO)

The measures of statistical significance presented above are not well suited to real-world applications, where the densities and sizes of the the dense components may vary significantly, and where we desire smoother transitions. Motivated by these shortcomings, we propose an alternate formulation that combines information about the size *and* density of the overlapping set, defined by the following probabilistic experiment:

- Consider a set  $V$  of vertices, with size  $n$ , and a distinguished subset  $\hat{B}$  of size  $|B|$ . Choose uniformly at random from  $V$  a subset  $\hat{A}$  of size  $|A|$ .
- In the set  $\hat{A} \cup \hat{B}$ , choose uniformly at random  $|e(A \cup B)|$  edges to insert.

Then, we define the combined density/overlap  $p$ -value to be the probability that the resulting overlap  $\hat{A} \cap \hat{B}$  has size at least that of the observed overlap  $|Z|$  *and* that the density of the overlapping set  $\delta(\hat{A} \cap \hat{B})$  is at least the observed overlap density  $\delta(Z)$ . That is, it is given by

$$p_{CoDO} = \Pr[|\hat{A} \cap \hat{B}| \geq |Z| \cap \delta(\hat{A} \cap \hat{B}) \geq \delta(Z)]$$

By conditioning on the size of the overlap, we can get an explicit formula for this  $p$ -value in terms of hypergeometric tails:

$$p_{CoDO} = \sum_{j=|Z|}^{\min\{|\hat{A}|, |\hat{B}|\}} \Pr[|\hat{A} \cap \hat{B}| = j] \cdot \Pr[\delta(\hat{A} \cap \hat{B}) \geq \delta(Z) | |\hat{A} \cap \hat{B}| = j] \quad (2)$$

The first factor in the terms of the sum is easily seen to be a hypergeometric probability mass:

$$|\hat{A} \cap \hat{B}| \sim \text{Hypergeometric}(n, |B|, |A|),$$

which is the number of successes when one draws without replacement  $|A|$  samples from a collection of  $n$  items,  $|B|$  of which are successes. Analogously, the second factor is a hypergeometric probability tail, since

$$e(\hat{A} \cap \hat{B}) \sim \text{Hypergeometric}\left(\binom{|A| + |B| - j}{2}, \binom{j}{2}, e(A \cup B)\right).$$

Thus, the sum is explicitly computable, and tail bounds can be applied if approximations are desired. Note that the expression (2) involves the observed sizes and densities of the constituent subgraphs, as well as the overlaps, through the distributions of  $|\hat{A} \cap \hat{B}|$  and  $\delta(\hat{A} \cap \hat{B})$ .

#### 3.3.1 Characteristics of the CoDO $p$ -Value

We investigate in more detail some interesting properties of the CoDO  $p$ -value defined above.

An elementary observation is that, whenever the density of overlap subgraph  $\delta(Z)$  is low enough, the conditional probabilities in the  $p_{CoDO}$  formula degenerate to 1, and it becomes equivalent to the hypergeometric tail  $p$ -value. Similarly, when the overlap size  $|Z|$  is small, all probabilities degenerate to 1.

Next, we look at the behavior of  $p_{CoDO}$  when we fix  $\delta(Z) = \rho_Z \in (0, 1)$  and vary the observed overlap size  $|Z|$ . We consider  $z = \frac{|Z|}{M} \in (0, 1)$ , where we define  $M = \min\{|A|, |B|\}$ .

We have the following lemma:

**Lemma 1** (Monotone decrease as  $|Z|$  increases). *For  $\rho_Z$  fixed as above, for  $z_1 < z_2$ , we have*

$$p_{CoDO}(z_1, \rho_Z) \geq p_{CoDO}(z_2, \rho_Z).$$

*Proof.* We have

$$\begin{aligned} p_{CoDO}(z_1, \rho_Z) - p_{CoDO}(z_2, \rho_Z) &= \sum_{j=z_1 M}^M \Pr[|\hat{A} \cap \hat{B}| = j] \Pr[\delta(\hat{A} \cap \hat{B}) \geq \rho_Z | |\hat{A} \cap \hat{B}| = j] \\ &\quad - \sum_{j=z_2 M}^M \Pr[|\hat{A} \cap \hat{B}| = j] \Pr[\delta(\hat{A} \cap \hat{B}) \geq \rho_Z | |\hat{A} \cap \hat{B}| = j] \\ &= \sum_{j=z_1 M}^{z_2 M-1} \Pr[|\hat{A} \cap \hat{B}| = j] \Pr[\delta(\hat{A} \cap \hat{B}) \geq \rho_Z | |\hat{A} \cap \hat{B}| = j] \end{aligned}$$

Since all of the terms of the final sum are at least 0, this completes the proof.  $\square$

It then becomes interesting to investigate the point at which  $p_{CoDO}$  transitions from insignificant to significant as a function of  $z$ . To do this we, define the *threshold point*  $\theta(\rho_Z)$  by

$$\theta(\rho_Z) = z_* = \inf\{z \in (0, 1) : p_{CoDO}(z, \rho_Z) \leq 1/2\}.$$

Note that, whenever it exists, the threshold point is unique, because of the monotonicity property just proven. In principle, the asymptotic value of the threshold point, as a function of  $\rho_Z$ , is explicitly computable. However, the fact that  $\mathbb{E}[e(\hat{A} \cap \hat{B}) | |\hat{A} \cap \hat{B}| = j]$  is an increasing function of  $j$  whenever  $j$  is large enough makes this somewhat subtle. The following upper bound can be established:

**Theorem 1.** *With  $\rho_Z \in (0, 1)$  fixed and*

$$M = \min\{|A|, |B|\} \gg \sqrt{n},$$

*we have, for any fixed  $\delta > 0$ ,*

$$\theta(\rho_Z) \leq \frac{\mathbb{E}[|\hat{A} \cap \hat{B}|]}{M} (1 + \delta).$$

To prove this, we will need the following lemma giving large deviation bounds for the hypergeometric distribution [16]. It is a consequence of the Hoeffding inequality [17].

**Lemma 2** (Hypergeometric large deviations). *Let  $X \sim \text{Hypergeometric}(N, K, n)$  (the parameters denote the population size, number of successes, and number of trials, respectively). Let  $p = K/N$ . Then we have the following tail bound:*

$$\Pr[X \geq (1 + \delta)\mathbb{E}[X]] \leq \exp(-n\kappa((1 + \delta)p, p)),$$

*where  $\delta$  is any number greater than 0, and  $\kappa(a, b)$  is given by (1).*

*of Theorem 1.* Let  $z = (1 + \delta) \frac{\mathbb{E}[|\hat{A} \cap \hat{B}|]}{M}$ . We have the following upper bound on  $p_{CoDO}(z, \rho_Z)$ .

$$p_{CoDO}(z, \rho_Z) \leq \sum_{j=zM}^M \Pr[|\hat{A} \cap \hat{B}| = j] p_* = \Pr[|\hat{A} \cap \hat{B}| \geq zM] p_*, \quad (3)$$

where we define

$$p_* = \max_j \left\{ \Pr[\delta(\hat{A} \cap \hat{B}) \geq \rho_Z | |\hat{A} \cap \hat{B}| = j] \right\},$$

where the maximum is taken over all terms  $j$  in the sum (this is allowed because all of the probabilities are non-negative).

Next, we apply Lemma 2 to further upper bound (3) by

$$p_{CoDO}(z, \rho_Z) \leq \Pr[|\hat{A} \cap \hat{B}| \geq (1 + \delta)\mathbb{E}[|\hat{A} \cap \hat{B}|]]p_* \leq \exp(-|A|\kappa((1 + \delta)|B|/n, |B|/n))p_*.$$

We now examine the role of the relative entropy factor in the exponent. As  $|B|/n \xrightarrow{n \rightarrow \infty} 0$ , its first term asymptotically dominates:

$$\kappa((1 + \delta)|B|/n, |B|/n) = (1 + \delta)|B|/n \log(1 + \delta) + o(1).$$

Using the fact that  $|B| \geq M \gg (\sqrt{n})$ , we then have that the entire exponent tends to  $-\infty$ . That is, provided that  $n$  is large enough,  $p_{CoDO}(z, \rho_Z)$  is arbitrarily small and, in particular, is less than  $1/2$ .  $\square$

Note that, in the proof of this theorem, we did not use any information about the conditional probability in the  $p_{CoDO}$  formula. This is another incarnation of the phenomenon that the conditional expectation of the number of edges in  $\hat{A} \cap \hat{B}$  increases with  $j$ . Thus, broadening the applicability of this bound requires a more careful study of how the conditional probabilities behave as  $j$  increases.

The behavior of  $p_{CoDO}(z, \rho_Z)$  for fixed and varying  $\rho_Z$  is also of interest. We have an analogue of Lemma 1:

**Lemma 3** (Monotone decrease as  $\rho_Z$  increases). *For fixed  $z \in (0, 1)$ , we have, for  $\rho_1 < \rho_2$ ,*

$$p_{CoDO}(z, \rho_1) \geq p_{CoDO}(z, \rho_2).$$

*Proof.* We have

$$\begin{aligned} & p_{CoDO}(z, \rho_1) - p_{CoDO}(z, \rho_2) \\ &= \sum_{j=zM}^M \Pr[|\hat{A} \cap \hat{B}| = j] \cdot \left[ \Pr[\delta(\hat{A} \cap \hat{B}) \geq \rho_1 | |\hat{A} \cap \hat{B}| = j] - \Pr[\delta(\hat{A} \cap \hat{B}) \geq \rho_2 | |\hat{A} \cap \hat{B}| = j] \right] \end{aligned}$$

By monotonicity of the CDF of a random variable, the difference of conditional probabilities is non-negative, so that the entire sum is positive. This completes the proof.  $\square$

As an analogue of the definitions for fixed  $\rho_Z$  and varying  $z$ , we can define a threshold point  $\phi(z)$ . We can also give analogous bounds for the location of the threshold point. The details of the derivation are entirely similar to those for fixed  $z$  and varying  $\rho_Z$ .

## 4 Experimental Validation

We provide a detailed experimental evaluation of CoDO, in comparison with other measures, in the context of synthetic datasets, datasets derived from social networks, and datasets from biomolecular interactions. We use synthetic datasets to characterize the dependence of our results of choices of parameters. We use the other datasets to demonstrate the robustness and application scope of our framework, even in cases where the networks do not follow an Erdos-Renyi model.

### 4.1 Assessing Overlap in Synthetic Graphs

To evaluate the performance of each measure (and associated method) on a controlled dataset, we create a synthetic test case that models the behavior of functional modules at a smaller scale. We sample an ER graph  $G_R$  from  $\mathcal{G}(n, p)$  with  $n = 80$  and  $p = \frac{3}{n}$ . We embed two modules with the same high density ( $\rho_A = \rho_B = 10p$ ) of size 50 and 40 vertices in  $G_R$  with varying levels of overlaps and density. For the overlap, the expected value of hypergeometric distribution for this setting is  $\frac{n_A n_B}{n} = 25$ , where  $n_A$  and  $n_B$  are the number of vertices in modules  $A$  and  $B$ , respectively. We chose overlap sizes of 20 and 30 vertices, which are below and above the expected value of number of overlapping vertices, respectively, to represent low and high overlaps. For simulating the density of overlap region, we tested a sparse overlap by setting  $\rho_Z = 2p = \frac{\rho_A}{5}$  and a dense overlap by setting  $\rho_Z = 10p = 2\rho_A$ . For each of these four settings, we computed  $HGT$ ,  $ERD$ , and  $CoDO$   $p$ -values, the results of which are illustrated in Figure 1. Each subfigure represents the adjacency matrix of the constructed test case, with modules  $A$  and  $B$  color-coded as blue and green, respectively. The

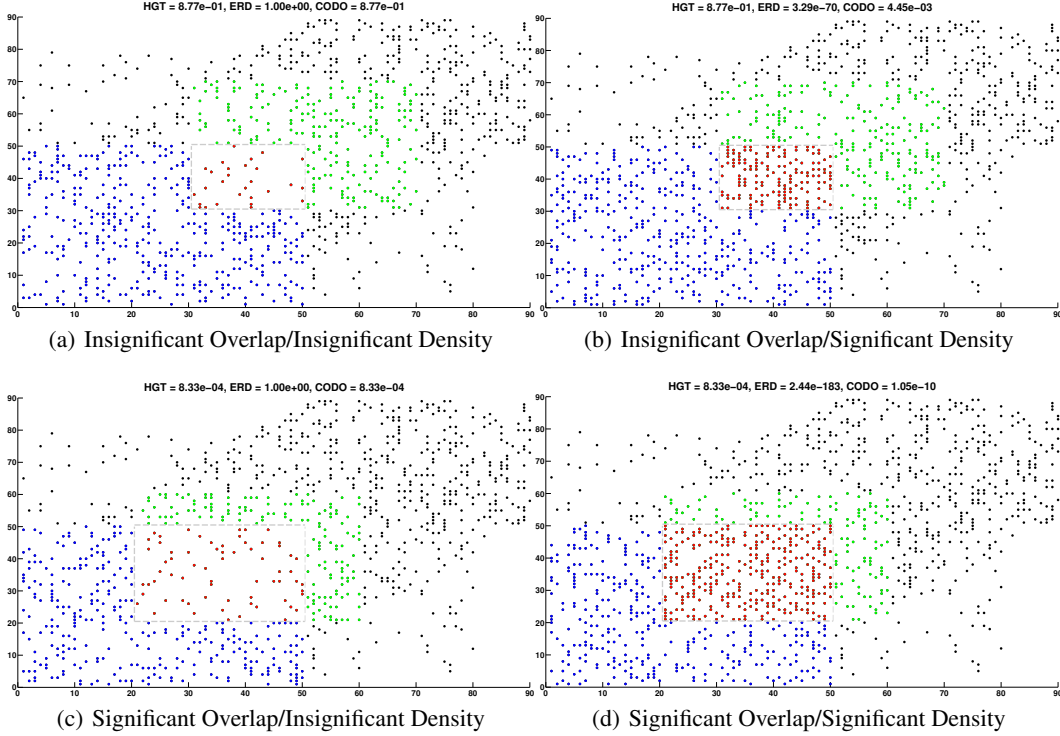


Figure 1: Synthetic graphs with varying levels of density/overlap. Illustrated here is the adjacency matrix of generated graphs. Black, green, and blue dots (edges) belong to the background, the first, and the second implanted subgraphs, respectively, whereas red dots represent edges in the overlap.

overlap set of the two modules is separately marked in red. The first notable observation in our experiments is that *ERD* has a sharp transition from  $p$ -value of 1 to 0 as we vary  $\rho_Z$  from  $\rho$  to  $2\rho$ . On the other hand, *CoDO* exhibits a smooth transition as we increase  $\rho_Z$ .

When neither the overlap, nor its density are significant, as in Figure 1(a), the  $p$ -values from all methods are close to 1 (as expected). When the overlap size is significant but density of overlap is not (Figure 1(c)), *HGT* and *CoDO* yield identical significant results, while *ERD* does not identify the overlap as significant ( $p$ -value approaching 1). On the other hand, when size of overlap is not significant but its density is significant, as in Figure 1(b), *HGT*  $p$ -value is close to one, whereas both *ERD* and *CoDO* identify the overlap as significant, the difference being that *ERD*  $p$ -value abruptly changes from 1 to  $3.29 \times 10^{-70}$ , but *CoDO* smoothly transitions from  $8.77 \times 10^{-1}$  to  $4.45 \times 10^{-3}$  demonstrating much better discriminating power. Finally, when both the size of the overlap and its density are significant, *CoDO*  $p$ -value is more significant than in both cases in Figures 1(b) and Figure 1(c), and it is more significant than *HGT*  $p$ -value alone.

## 4.2 Significance of Overlap Among Social Circles

Users in social networks are connected not only to their close friends, but also family members, schoolmates, and colleagues, among others. Organizing friends into communities, or *social circles*, is one of the common approaches to organizing, predicting, and recommending contacts. A user's circles are typically *overlapping*, and a user can belong to one, many, or none of the circles. Given a set of overlapping circles, we are interested in assessing whether there is a significant interdependency among the circles, or that the observed overlap is simply due to the existence of *bridge nodes* or *party alters*.

To evaluate our method, we use three sets of social networks derived from Facebook, Google+, and Twitter[3]. Each of these networks is centered around a focal user, or *ego*, along with all other users it directly connects to. This collection of neighbor nodes is also known as the set of *alters*. The induced subgraph of friendships among *alter nodes*, is known as an *ego net*. The Facebook dataset is *fully*



Table 1: Statistics of ego nets

Dataset	# Nets	Avg # Users	Avg # Circles	# Over. Circles	# Signif. Overlaps
Facebook	10	416.7	19.3	415	92
Google+	124	1,945.7	3.5	1,170	437
Twitter	834	137.6	4.2	14,610	5,055

*observed*, in the sense that each user was asked to manually identify all coherent groups among their friends. We refer to these coherent groups of friends as *circles*. In contrast to the Facebook dataset, circles in Google+ and Twitter are restricted to publicly visible and explicitly designated circles for each ego net.

To establish a ground truth for interdependency among pairs of circles, we use common features of users to assign a  $p$ -value to their feature overlap. For a given ego net  $G = (V, E)$  together with a  $K$ -dimensional feature vector for each node, we compute the total number of feature pairs between any pair of alters as  $\pi = K \binom{|V|}{2}$ . Moreover, we can compute the total number of features common to any two alters as:

$$\xi = \sum_{\substack{i, j \in V \\ i < j}} \sum_{f \in \mathcal{F}} \mathbb{1}_f(i, j) \quad (4)$$

where  $V$  is the set of alter vertices,  $\mathcal{F}$  is the set of features, and  $\mathbb{1}_f(i, j)$  is an indicator function that is one, if alter  $i$  and alter  $j$  share feature  $f$ , and zero, otherwise. Given an overlap set of  $O$ , we can similarly define  $\pi_O = K \binom{|O|}{2}$  and:

$$\xi_O = \sum_{\substack{i, j \in O \\ i < j}} \sum_{f \in \mathcal{F}} \mathbb{1}_f(i, j) \quad (5)$$

Using this notion, we can define the significance of feature overlap among two circles as:

$$p - val(O) = \sum_{x=\xi_O}^{\min(\xi, \pi_O)} \frac{\binom{\xi}{x} \binom{\pi - \xi}{\pi_O - x}}{\binom{\pi}{\pi_O}} \quad (6)$$

which is simply the tail of the Hypergeometric distribution, measuring the probability of observing  $\xi_O$  or more common features in a random set of size  $O$  (with  $\pi_O$  feature pairs), if we were sampling at random without replacement from all possible feature pairs. For each dataset, we compute the significance of all circle pairs for every ego net and correct for multiple hypothesis testing using *Bonferroni* method. Finally, we define a pair of circles as *significantly correlated* if  $p$ -value of their overlap is smaller than or equal to the threshold 0.05. Table 1 summarizes the statistics of these networks. Here, we remove all ego nets with less than two circles or those without at least one common feature between alters. Among the three datasets, Facebook had the lowest percentage of significant overlaps ( $< 25\%$ ) compared to the other two datasets (35% and 37%). This may be attributed to method used for identifying circles in the Facebook dataset.

Next, to evaluate *HGT*, *ERD*, and *CoDO* methods, we assess all pairs of overlapping circles and sort them based on their significance. For pairs with similar  $p$ -value, we randomly order them. Using the gold-standard set of *significantly correlated* circles, computed using feature vector overlaps, we compute the *receiver operating characteristic (ROC)* for each dataset individually. These are presented in Figure 2. Each method is annotated with its area under the curve (AUC) to simplify comparison. In all three datasets, *CoDO* outperforms the other two methods. It is worth noting that *CoDO* is designed to capture signals from both network density and the overlap size. To illustrate this point, we emphasize the behavior of *CoDO* over the Facebook dataset as an example. For small FPR values, density has higher signal than overlap (observed as the superior performance of ERD compared to HGT). In this regime, *CoDO* captures this signal and *outperforms* of *ERD*. When *ERD* signal plateaus ( $0.1 < FPR$ ), *CoDO* leverages network structure to outperform (*HGT*).

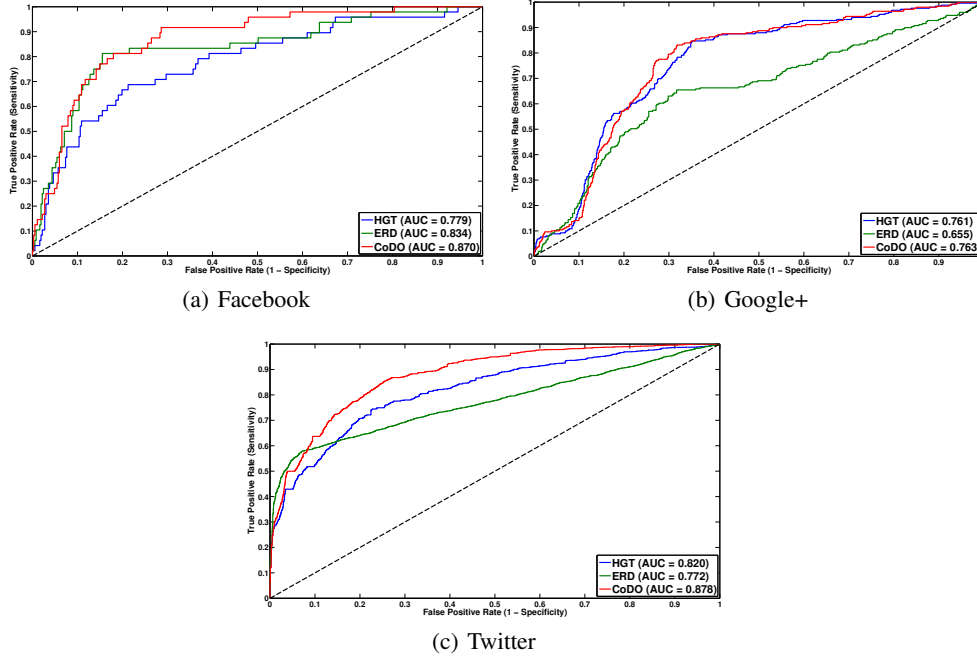


Figure 2: Overlapping circles in social networks predict common features

### 4.3 Crosstalk Among Biological Pathways

Biological pathways are higher order constructs that perform key cellular functions. These pathways are known to work in concert and their cross-talk regulates the systems level behavior of cells. To evaluate the interaction between different pathways and their shared mechanisms, we seek to construct a pathway-pathway interaction map that represents the extent of overlap among different pathways.

We construct a comprehensive human interactome from a recently published dataset by Han *et al.*[18], resulting in a network of 164,826 interactions among 8,872 proteins. We downloaded the set of KEGG pathways from MSigDB[19] and filtered pathways that have less than 10 corresponding vertices in the human interactome. The final dataset consists of 186 pathways with an average of  $\sim 66$  vertices per pathway. To evaluate the computed  $p$ -values, we use *co-transcriptional activity of pathway pairs* as a proxy for their functional relatedness. We downloaded the and processed tissue-specific RNASeq dataset from the Genotype-Tissue Expression (GTEx) project [20], which contains 2,916 samples from 30 different tissues/ cell types. We summarize tissue-specific expression of each pathway by averaging the expression of all its member genes in each tissue. Then, we define co-transcriptional activity of pathways by computing the Pearson’s correlation of these tissue-specific pathway expression signatures.

To evaluate different methods, we first compute the nonparametric Spearman’s correlation between the co-transcriptional activity scores and computed  $p$ -values in each method. This yields 0.28, 0.16, and 0.41 correlation scores for *HGT*, *ERD*, and *CoDO* methods, respectively, which shows that *CoDO* significantly outperforms the other two methods with respect to the co-transcriptional activity of significant pathway pairs. Next, to put the computed  $p$ -values by *CoDO* in context, we construct a pathway-pathway overlap network by thresholding pairwise overlaps at a stringent cutoff ( $p$ -value  $\leq 10^{-30}$ ) and use all pairs with significant  $p$ -values as edges in the network. This results in a network of 129 nodes (representing pathways) and 877 interactions (inferred as the significance of the pathway overlaps) among them. This network is shown in Figure 3. We use Cytoscape [21] to visualize the graph and MCODE [22] to cluster the network. Each cluster is color-coded independently, and four major clusters are manually annotated with the dominant function in each group. We observe that the set of top-ranked pathway pairs cluster together to form coherent groups with highly coherent functions. These significant overlaps reveal interesting functional connections, which can be used to understand pathology and identify novel drug targets.

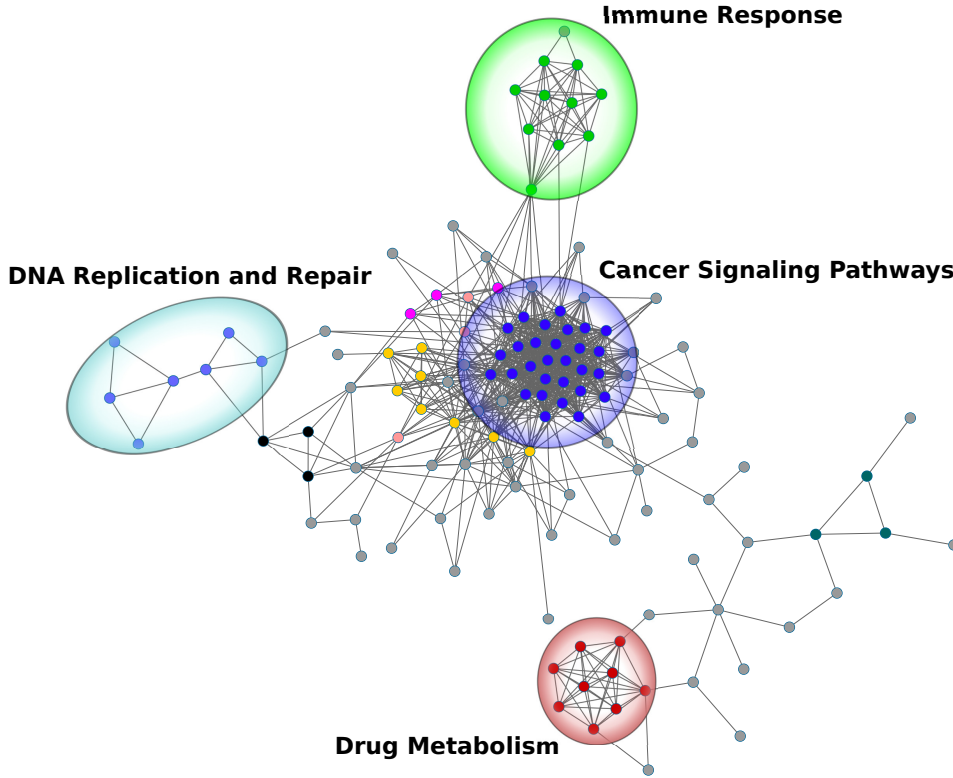


Figure 3: Pathway-pathway overlap graph

## 5 Conclusions

Assessing the statistical significance of observed overlaps between clusters in networks is an important substrate of many graph analytics problems. In this paper, we present detailed analysis and derivation of a  $p$ -value based measure of significance of cluster overlap. Unlike previous measures, our measure accounts for cluster sizes, densities, and density of overlap – all critical parameters of the problem. We show that our measure provides excellent discrimination, smooth transition, and robustness to wide parameter ranges. Using both synthetic and real datasets, we validate excellent performance of our analytical formulation.

Our work opens a number of avenues for continued explorations. These include formulations for alternate graph models, algorithms aimed at explicitly maximizing statistical significance of overlaps, and application validation in other contexts.

## References

- [1] M. Koyutürk, W. Szpankowski, and A. Grama, “Assessing significance of connectivity and conservation in protein interaction networks,” *Journal of Computational Biology*, vol. 14, no. 6, pp. 747–764, 2007.
- [2] S. Mohammadi et al., “Scope and limitations of yeast as a model organism for studying human tissue-specific pathways,” *BMC Systems Biology*, vol. 9, no. 1, p. 96, Dec. 2015.
- [3] J. Mcauley and J. Leskovec, “Discovering social circles in ego networks,” *ACM Trans. Knowl. Discov. Data*, vol. 8, no. 1, 4:1–4:28, Feb. 2014.
- [4] J. McDonald, *Handbook of Biological Statistics*. Baltimore, MD, USA: Sparky House Publishing, 2014.
- [5] V. E. Lee et al., “A survey of algorithms for dense subgraph discovery,” English, in *Managing and Mining Graph Data*, C. C. Aggarwal and H. Wang, Eds., vol. 40, ser. Advances in Database Systems, Springer US, 2010, pp. 303–336.
- [6] R. Arratia and E. S. Lander, “The distribution of clusters in random graphs,” *Advances in Applied Mathematics*, vol. 11, pp. 36–48, 1990.
- [7] B. Bollobás, *Random Graphs*, 2nd. Cambridge Studies in Advanced Mathematics, 2001.
- [8] D. Gamarnik and M. Sudan, “Limits of local algorithms over sparse random graphs,” *Proceedings of the 5th Conference on Innovations in Theoretical Computer Science*, pp. 369–376, 2014.
- [9] Y. Sun et al., “Crosstalk analysis of pathways in breast cancer using a network model based on overlapping differentially expressed genes,” *Experimental and Therapeutic Medicine*, May 2015.
- [10] Y. Wang et al., “Pathway crosstalk analysis of high-metastasis lung cancer cells,” *Tumori*,
- [11] S. Mohammadi, G. Kollias, and A. Grama, “Role of synthetic genetic interactions in understanding functional interactions among pathways,” in *Pacific Symposium on Biocomputing (PSB) 2012*, Jan. 2012, pp. 43–54.
- [12] A. N. Tegge, N. Sharp, and T. M. Murali, “Xtalk: a path-based approach for identifying crosstalk between signaling pathways,” *Bioinformatics*, btv549, Sep. 2015.
- [13] C. Mitrea et al., “Methods and approaches in the topology-based analysis of biological pathways,” *Frontiers in Physiology*, vol. 4, 2013.
- [14] X. Dong et al., “LEGO: a novel method for gene set over-representation analysis by incorporating network-based gene weights,” *Scientific Reports*, vol. 6, p. 18 871, Jan. 2016.
- [15] N. Alon and J. H. Spencer, *The probabilistic method*. Wiley.
- [16] V. Chvátal, “The tail of the hypergeometric distribution,” *Discrete Mathematics*, vol. 25, pp. 285–287, 1979.
- [17] A. Dembo and O. Zeitouni, *Large Deviations Techniques and Applications*, 2nd. Springer-Verlag, 2010.
- [18] J. Han et al., “ESEA: Discovering the Dysregulated Pathways based on Edge Set Enrichment Analysis,” *Scientific reports*, vol. 5, p. 13 044, 2015.
- [19] A. Subramanian et al., “Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles,” *Proceedings of the National Academy of Sciences*, vol. 102, no. 43, pp. 15 545–15 550, Oct. 2005.
- [20] K. G. Ardlie et al., “The Genotype-Tissue Expression (GTEx) pilot analysis: Multitissue gene regulation in humans,” *Science*, vol. 348, no. 6235, pp. 648–660, May 2015.
- [21] P. Shannon et al., “Cytoscape: A software Environment for integrated models of biomolecular interaction networks,” *Genome Research*, vol. 13, no. 11, pp. 2498–2504, 2003.
- [22] G. Bader and C. Hogue, “An automated method for finding molecular complexes in large protein interaction networks,” *BMC bioinformatics*, vol. 27, pp. 1–27, 2003.